```
In [1]:  # Initialize autograder
         # If you see an error message, you'll need to do
         # pip3 install otter-grader
         import otter
         grader = otter.Notebook()
```

# Project 3: Predicting Taxi Ride Duration

## Due Date: Wednesday 3/4/20, 11:59PM

**Collaboration Policy**

Data science is a collaborative activity. While you may talk with others about the project, we ask that you **write your solutions individually**. If you do discuss the assignments with others please **include their names** at the top of your notebook.

**Collaborators**: *list collaborators here*

# Score Breakdown

| Question | Points |
| --- | --- |
| 1b | 2 |
| 1c | 3 |
| 1d | 2 |
| 2a | 1 |
| 2b | 2 |
| 3a | 2 |
| 3b | 1 |
| 3c | 2 |
| 3d | 2 |
| 4a | 2 |
| 4b | 2 |
| 4c | 2 |
| 4d | 2 |
| 4e | 2 |
| 4f | 2 |
| 4g | 4 |
| 5b | 7 |
| 5c | 3 |
| Total | 43 |

# This Assignment

In this project, you will use what you've learned in class to create a regression model that predicts the travel time of a taxi ride in New York. Some questions in this project are more substantial than those of past projects.

After this project, you should feel comfortable with the following:

- The data science lifecycle: data selection and cleaning, EDA, feature engineering, and model selection.
- Using `sklearn` to process data and fit linear regression models.
- Embedding linear regression as a component in a more complex model.

First, let's import:

```
import numpy as np
import pandas as pd

import matplotlib.pyplot as plt
%matplotlib inline

import seaborn as sns
```

# The Data

Attributes of all yellow taxi (https://en.wikipedia.org/wiki/Taxicabs_of_New_York_City) trips in January 2016 are published by the NYC Taxi and Limosine Commission (https://www1.nyc.gov/site/tlc/about/tlc-trip-record-data.page).

The full data set takes a long time to download directly, so we've placed a simple random sample of the data into `taxi.db`, a SQLite database. You can view the code used to generate this sample in the `taxi_sample.ipynb` file included with this project (not required).

Columns of the `taxi` table in `taxi.db` include:

- `pickup_datetime`: date and time when the meter was engaged
- `dropoff_datetime`: date and time when the meter was disengaged
- `pickup_lon`: the longitude where the meter was engaged
- `pickup_lat`: the latitude where the meter was engaged
- `dropoff_lon`: the longitude where the meter was disengaged
- `dropoff_lat`: the latitude where the meter was disengaged
- `passengers`: the number of passengers in the vehicle (driver entered value)
- `distance`: trip distance
- `duration`: duration of the trip in seconds

Your goal will be to predict `duration` from the pick-up time, pick-up and drop-off locations, and distance.

# Part 1: Data Selection and Cleaning

In this part, you will limit the data to trips that began and ended on Manhattan Island (map (https://www.google.com/maps/place/Manhattan,+New+York,+NY/@40.7590402,-74.0394431,12z/data=!3m1!4b1!73.9712488)).

The below cell uses a SQL query to load the `taxi` table from `taxi.db` into a Pandas DataFrame called `all_taxi`.

It only includes trips that have **both** pick-up and drop-off locations within the boundaries of New York City:

- Longitude is between -74.03 and -73.75 (inclusive of both boundaries)
- Latitude is between 40.6 and 40.88 (inclusive of both boundaries)

You don't have to change anything, just run this cell.

```python
In [3]: import sqlite3

        conn = sqlite3.connect('taxi.db')
        lon_bounds = [-74.03, -73.75]
        lat_bounds = [40.6, 40.88]

        c = conn.cursor()

        my_string = 'SELECT * FROM taxi WHERE'

        for word in ['pickup_lat', 'AND dropoff_lat']:
            my_string += ' {} BETWEEN {} AND {}'.format(word, lat_bounds[0], lat_
        bounds[1])

        for word in ['AND pickup_lon', 'AND dropoff_lon']:
            my_string += ' {} BETWEEN {} AND {}'.format(word, lon_bounds[0], lon_
        bounds[1])

        c.execute(my_string)

        results = c.fetchall()

        row_res = conn.execute('select * from taxi')
        names = list(map(lambda x: x[0], row_res.description))


        all_taxi = pd.DataFrame(results)
        all_taxi.columns = names
        all_taxi.head()
```
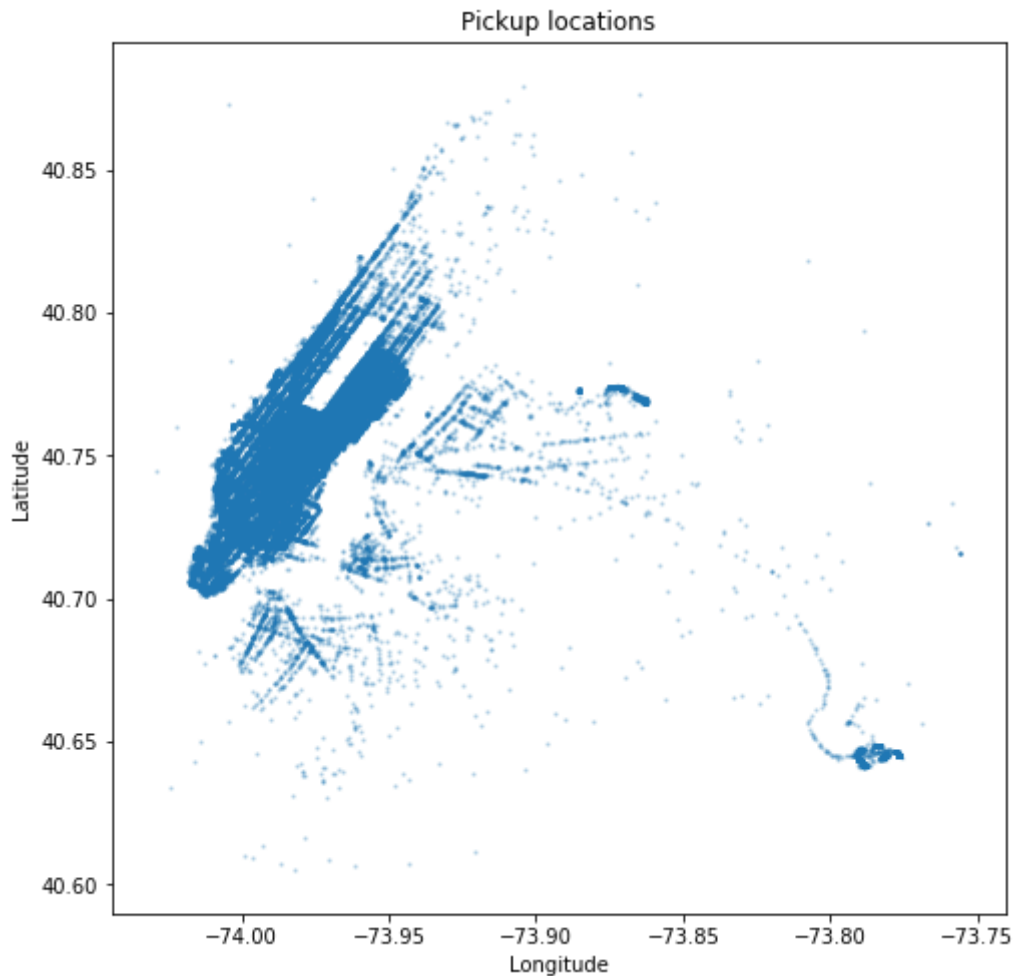
Out[3]:

| | pickup_datetime | dropoff_datetime | pickup_lon | pickup_lat | dropoff_lon | dropoff_lat | passenger |
|---|---|---|---|---|---|---|---|
| 0 | 2016-01-30 22:47:32 | 2016-01-30 23:03:53 | -73.988251 | 40.743542 | -74.015251 | 40.709808 | |
| 1 | 2016-01-04 04:30:48 | 2016-01-04 04:36:08 | -73.995888 | 40.760010 | -73.975388 | 40.782200 | |
| 2 | 2016-01-07 21:52:24 | 2016-01-07 21:57:23 | -73.990440 | 40.730469 | -73.985542 | 40.738510 | |
| 3 | 2016-01-01 04:13:41 | 2016-01-01 04:19:24 | -73.944725 | 40.714539 | -73.955421 | 40.719173 | |
| 4 | 2016-01-08 18:46:10 | 2016-01-08 18:54:00 | -74.004494 | 40.706989 | -74.010155 | 40.716751 | |

A scatter plot of pickup locations shows that most of them are on the island of Manhattan. The empty white rectangle is Central Park; cars are not allowed there.

```
In [4]:  def pickup_scatter(t):
             plt.scatter(t['pickup_lon'], t['pickup_lat'], s=2, alpha=0.2)
             plt.xlabel('Longitude')
             plt.ylabel('Latitude')
             plt.title('Pickup locations')

         plt.figure(figsize=(8, 8))
         pickup_scatter(all_taxi)
```



The two small blobs outside of Manhattan with very high concentrations of taxi pick-ups are airports.

## Question 1b

Create a DataFrame called `clean_taxi` that only includes trips with a positive passenger count, a positive distance, a duration of at least 1 minute and at most 1 hour, and an average speed of at most 100 miles per hour. Inequalities should not be strict (e.g., `<=` instead of `<`) unless comparing to 0.

*The provided tests check that you have constructed* `clean_taxi` *correctly.*

```
In [5]: clean_taxi = all_taxi[(all_taxi["passengers"] > 0) & (all_taxi["distance"
        ] > 0)
                            & (all_taxi["duration"] >= 60) & (all_taxi["duratio
        n"] <= 3600)
                            & (all_taxi["distance"]/(all_taxi["duration"] / 360
        0) <= 100)
                            ]
```

```
In [6]: grader.check("q1b")
```

Out[6]:  All tests passed!


## Question 1c (challenging)

Create a DataFrame called `manhattan_taxi` that only includes trips from `clean_taxi` that start and end within a polygon that defines the boundaries of Manhattan Island (https://www.google.com/maps/place/Manhattan,+New+York,+NY/@40.7590402,-74.0394431,12z/data=!3m1!4b1!73.9712488).

The vertices of this polygon are defined in `manhattan.csv` as (latitude, longitude) pairs, which are published here (https://gist.github.com/baygross/5430626).

An efficient way to test if a point is contained within a polygon is described on this page (http://alienryderflex.com/polygon/). There are even implementations on that page (though not in Python). Even with an efficient approach, the process of checking each point can take several minutes. It's best to test your work on a small sample of `clean_taxi` before processing the whole thing. (To check if your code is working, draw a scatter diagram of the (lon, lat) pairs of the result; the scatter diagram should have the shape of Manhattan.)

*The provided tests check that you have constructed `manhattan_taxi` correctly. It's not required that you implement the `in_manhattan` helper function, but that's recommended. If you cannot solve this problem, you can still continue with the project; see the instructions below the answer cell.*

```
In [7]:  polygon = pd.read_csv('manhattan.csv')
         polyCorners = len(polygon) # number of corners
         polyX = np.array(polygon["lon"]) # horizontal coordinates of corners
         polyY = np.array(polygon["lat"]) # vertical coordinates of corners
         constant = np.empty(polyCorners) # precalculated constants
         multiple = np.empty(polyCorners) # precalculated multipliers

         def precalc_values():
             j = polyCorners-1
             for i in range(polyCorners):
                 if polyY[j] == polyY[i]:
                     constant[i] = polyX[i]
                     multiple[i] = 0
                 else:
                     constant[i] = polyX[i] - (polyY[i] * polyX[j]) / (polyY[j] -
         polyY[i]) \
                                      + (polyY[i] * polyX[i]) / (polyY[j] - polyY[i
         ])
                     multiple[i] = (polyX[j] - polyX[i]) / (polyY[j] - polyY[i])
                 j = i

         def in_manhattan(x, y):
             """Whether a longitude-latitude (x, y) pair is in the Manhattan polyg
         on."""
             j = polyCorners - 1
             oddNodes = False
             for i in range(polyCorners):
                 if (polyY[i] < y and polyY[j] >= y) or (polyY[j] < y and polyY[i]
          >= y):
                     oddNodes ^= (y * multiple[i] + constant[i] < x)
                 j = i
             return oddNodes

         # Recommended: Then, apply this function to every trip to filter clean_ta
         xi.
         precalc_values()
         pickup_X = clean_taxi["pickup_lon"].to_numpy()
         pickup_Y = clean_taxi["pickup_lat"].to_numpy()
         dropoff_X = clean_taxi["dropoff_lon"].to_numpy()
         dropoff_Y = clean_taxi["dropoff_lat"].to_numpy()
         inManhattan = np.zeros(pickup_X.shape, dtype=bool)

         coords = np.column_stack((pickup_X, pickup_Y, dropoff_X, dropoff_Y))
         for i, c in enumerate(coords):
             if in_manhattan(c[0], c[1]) and in_manhattan(c[2], c[3]):
                 inManhattan[i] = True

         manhattan_taxi = clean_taxi[inManhattan]
```

```
In [8]:  grader.check("q1c")
```
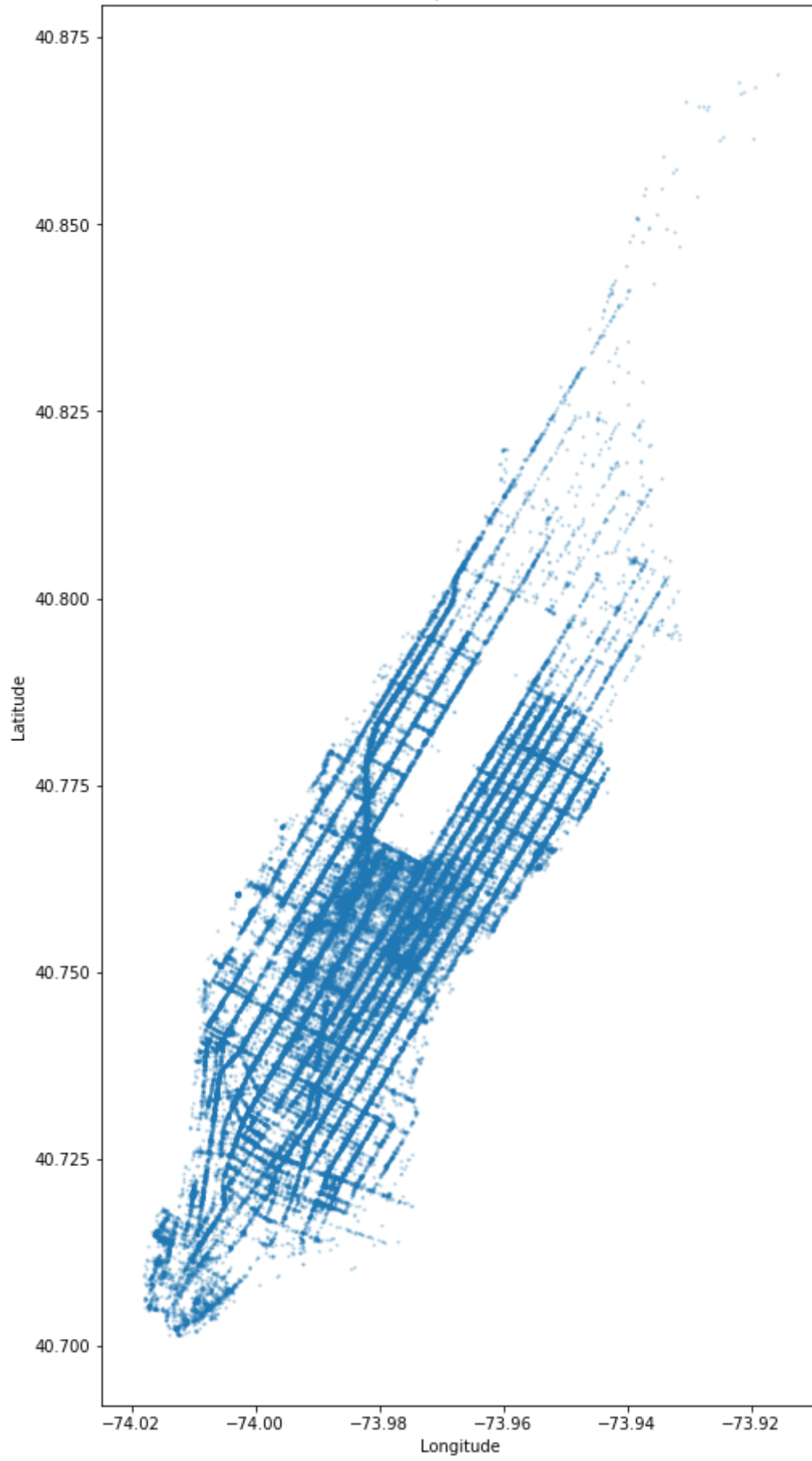
Out[8]:  **All tests passed!**

If you are unable to solve the problem above, have trouble with the tests, or want to work on the rest of the project before solving it, run the following cell to load the cleaned Manhattan data directly. (Note that you may not solve the previous problem just by loading this data file; you have to actually write the code.)

```
In [9]:  # manhattan_taxi = pd.read_csv('manhattan_taxi.csv')
```

A scatter diagram of only Manhattan taxi rides has the familiar shape of Manhattan Island.

```
In [10]: plt.figure(figsize=(8, 16))
         pickup_scatter(manhattan_taxi)
```

Pickup locations

# Question 1d

Print a summary of the data selection and cleaning you performed. **Your Python code should not include any number literals, but instead should refer to the shape of `all_taxi`, `clean_taxi`, and `manhattan_taxi`.**

E.g., you should print something like: "Of the original 1000 trips, 21 anomalous trips (2.1%) were removed through data cleaning, and then the 600 trips within Manhattan were selected for further analysis."

(Note that the numbers in the example above are not accurate.)

One way to do this is with Python's f-strings. For instance,

```
name = "Joshua"
print(f"Hi {name}, how are you?")
```

prints out `Hi Joshua, how are you?`.

**Please ensure that your Python code does not contain any very long lines, or we can't grade it.**

*Your response will be scored based on whether you generate an accurate description and do not include any number literals in your Python expression, but instead refer to the dataframes you have created.*

```
In [11]:  n_all = all_taxi.shape[0]
          n_cleaned = clean_taxi.shape[0]
          p_cleaned = (1-n_cleaned/n_all)*100
          n_filtered = manhattan_taxi.shape[0]
          p_filtered = (n_filtered/n_all)*100
          print(f"Of the original {n_all} rows of data, {n_cleaned} rows were selec
          ted (or {p_cleaned:.1f}% removed) according to the constraints detailed i
          n 1b.")
          print(f"Of the cleaned data, {n_filtered} rows were selected for analysi
          s, corresponding to trips that take place entirely within Manhattan.")
          print(f"This represents {p_filtered:.2f}% of the original data.")
```

```
Of the original 97692 rows of data, 96445 rows were selected (or 1.3% r
emoved) according to the constraints detailed in 1b.
Of the cleaned data, 82800 rows were selected for analysis, correspondi
ng to trips that take place entirely within Manhattan.
This represents 84.76% of the original data.
```

# Part 2: Exploratory Data Analysis

In this part, you'll choose which days to include as training data in your regression model.

Your goal is to develop a general model that could potentially be used for future taxi rides. There is no guarantee that future distributions will resemble observed distributions, but some effort to limit training data to typical examples can help ensure that the training data are representative of future observations.

January 2016 had some atypical days. New Year's Day (January 1) fell on a Friday. MLK Day was on Monday, January 18. A historic blizzard (https://en.wikipedia.org/wiki/January_2016_United_States_blizzard) passed through New York that month. Using this dataset to train a general regression model for taxi trip times must account for these unusual phenomena, and one way to account for them is to remove atypical days from the training data.

## Question 2a

Add a column labeled `date` to `manhattan_taxi` that contains the date (but not the time) of pickup, formatted as a `datetime.date` value (docs (https://docs.python.org/3/library/datetime.html#date-objects)).

*The provided tests check that you have extended* `manhattan_taxi` *correctly.*

```
In [12]:   from datetime import datetime
           # def getDate(dt):
           #     return datetime.strptime(dt, '%Y-%m-%d %H:%M:%S').date()
           # manhattan_taxi["date"] = manhattan_taxi["pickup_datetime"].map(getDate)
           manhattan_taxi["date"] = pd.to_datetime(manhattan_taxi["pickup_datetime"
           ]).dt.date
           manhattan_taxi.head()
```

```
C:\Users\LinusLam\Anaconda3\lib\site-packages\ipykernel_launcher.py:5:
SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: http://pandas.pydata.org/pandas-d
ocs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
  """
```

Out[12]:

| | pickup_datetime | dropoff_datetime | pickup_lon | pickup_lat | dropoff_lon | dropoff_lat | passengers |
|---|---|---|---|---|---|---|---|
| 0 | 2016-01-30 22:47:32 | 2016-01-30 23:03:53 | -73.988251 | 40.743542 | -74.015251 | 40.709808 | |
| 1 | 2016-01-04 04:30:48 | 2016-01-04 04:36:08 | -73.995888 | 40.760010 | -73.975388 | 40.782200 | |
| 2 | 2016-01-07 21:52:24 | 2016-01-07 21:57:23 | -73.990440 | 40.730469 | -73.985542 | 40.738510 | |
| 4 | 2016-01-08 18:46:10 | 2016-01-08 18:54:00 | -74.004494 | 40.706989 | -74.010155 | 40.716751 | |
| 5 | 2016-01-02 12:39:57 | 2016-01-02 12:53:29 | -73.958214 | 40.760525 | -73.983360 | 40.760406 | |

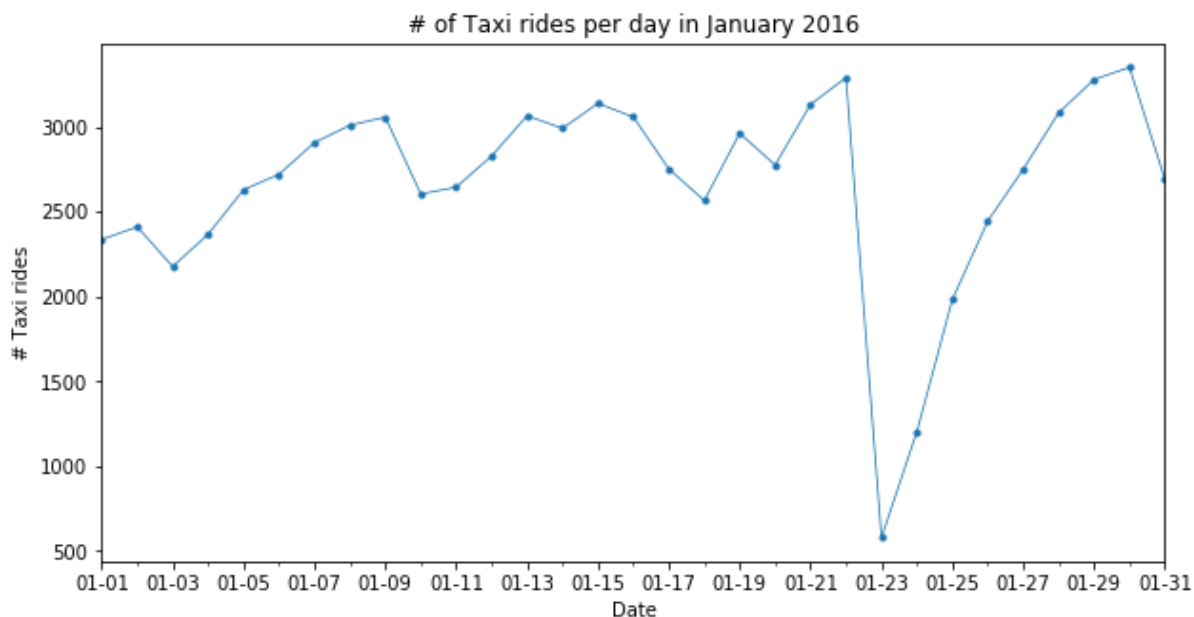```
In [13]:   grader.check("q2a")
```

Out[13]:   All tests passed!

## Question 2b

Create a data visualization that allows you to identify which dates were affected by the historic blizzard of January 2016. Make sure that the visualization type is appropriate for the visualized data.

As a hint, consider how taxi usage might change on a day with a blizzard. How could you visualize/plot this?

```
In [14]: date_freq = manhattan_taxi.groupby(["date"]).size()
         ax = date_freq.plot(figsize=(10,5), style=".-", lw=0.75)
         ax.set_xlabel("Date")
         ax.set_ylabel("# Taxi rides")
         ax.set_title("# of Taxi rides per day in January 2016")
         import matplotlib.dates as mdates
         ax.xaxis.set_major_locator(mdates.DayLocator(interval=2))
         ax.xaxis.set_major_formatter(mdates.DateFormatter('%m-%d'))
         ax.xaxis.set_minor_locator(mdates.DayLocator(interval=1))
```



Finally, we have generated a list of dates that should have a fairly typical distribution of taxi rides, which excludes holidays and blizzards. The cell below assigns `final_taxi` to the subset of `manhattan_taxi` that is on these days. (No changes are needed; just run this cell.)

```
In [15]:   import calendar
           import re

           from datetime import date

           atypical = [1, 2, 3, 18, 23, 24, 25, 26]
           typical_dates = [date(2016, 1, n) for n in range(1, 32) if n not in atypi
           cal]
           typical_dates

           print('Typical dates:\n')
           pat = '   [1-3]|18 | 23| 24|25 |26 '
           print(re.sub(pat, '   ', calendar.month(2016, 1)))

           final_taxi = manhattan_taxi[manhattan_taxi['date'].isin(typical_dates)]

           Typical dates:

                January 2016
           Mo Tu We Th Fr Sa Su

            4  5  6  7  8  9 10
           11 12 13 14 15 16 17
              19 20 21 22
                    27 28 29 30 31
```

You are welcome to perform more exploratory data analysis, but your work will not be scored. Here's a blank cell to use if you wish. In practice, further exploration would be warranted at this point, but the project is already pretty long.

```
In [16]:   # Optional: More EDA here
```

# Part 3: Feature Engineering

In this part, you'll create a design matrix (i.e., feature matrix) for your linear regression model. This is analagous to the pipelines you've built already in class: you'll be adding features, removing labels, and scaling among other things.

You decide to predict trip duration from the following inputs: start location, end location, trip distance, time of day, and day of the week (*Monday, Tuesday, etc.*).

You will ensure that the process of transforming observations into a design matrix is expressed as a Python function called `design_matrix`, so that it's easy to make predictions for different samples in later parts of the project.

Because you are going to look at the data in detail in order to define features, it's best to split the data into training and test sets now, then only inspect the training set.

```
In [17]:  import sklearn.model_selection

          train, test = sklearn.model_selection.train_test_split(
              final_taxi, train_size=0.8, test_size=0.2, random_state=42)
          print('Train:', train.shape, 'Test:', test.shape)

          Train: (53680, 10) Test: (13421, 10)
```
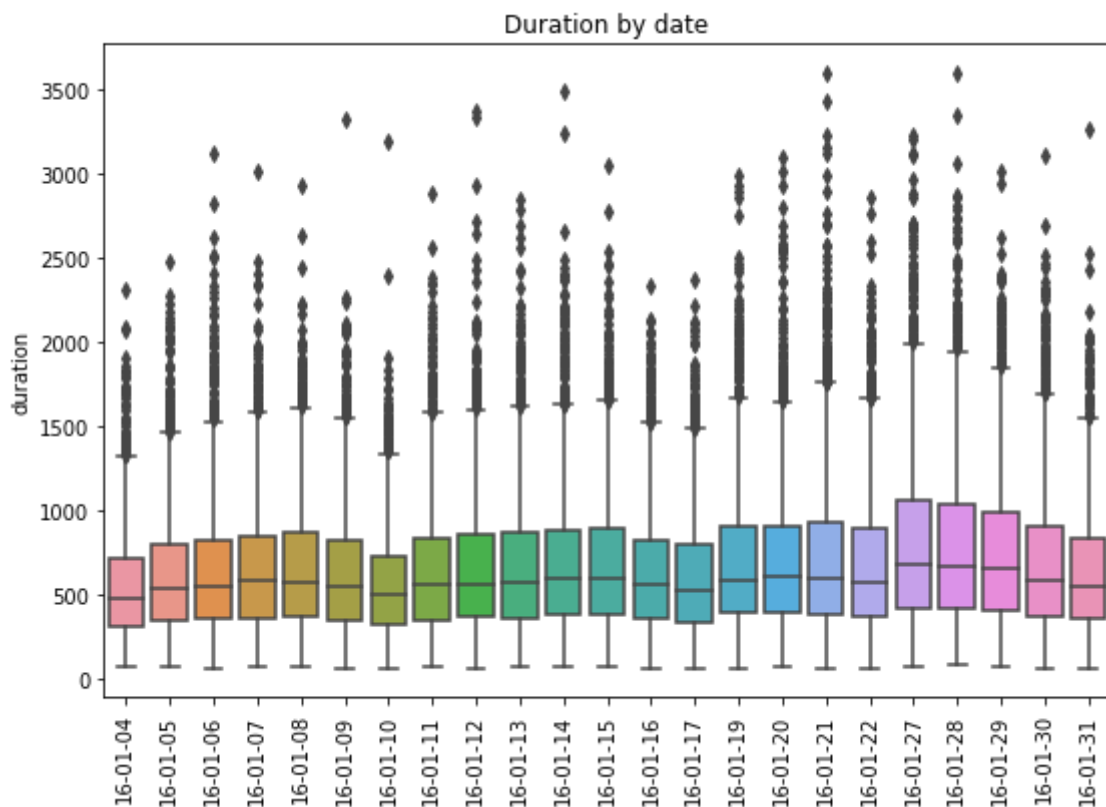
## Question 3a

Create a box plot that compares the distributions of taxi trip durations for each day **using `train` only**.
Individual dates shoud appear on the horizontal axis, and duration values should appear on the vertical axis.
Your plot should look like the one below.

You can generate this type of plot using `sns.boxplot`

```
In [18]:  plt.figure(figsize=(9,6))
          plt.xticks(rotation=90)
          dates = train["date"].sort_values(ascending=True).unique()
          ax = sns.boxplot(x=train["date"], y=train["duration"], order=dates)
          dformat = np.vectorize(lambda d : datetime.strftime(d, "%y-%m-%d"))
          ax.set_title("Duration by date")
          ax.set_xlabel(None)
          ax.set_xticklabels(dformat(dates))
          plt.show()
```

# Question 3b

In one or two sentences, describe the assocation between the day of the week and the duration of a taxi trip. Your answer should be supported by your boxplot above.

*Note*: The end of Part 2 showed a calendar for these dates and their corresponding days of the week.

The duration as well as the variance in duration tends to increase until the weekend, then return to their lowest levels by Monday.

In the two weeks from 1/4-1/17, duration peaked on Fridays. In the last two weeks of the month, duration peaked on Thursday and Wednesday respectively; however, it is worth noting that the Jan 2016 blizzard hit between these two weeks, and probably changed duration patterns immediately before/after.

Below, the provided `augment` function adds various columns to a taxi ride dataframe.

- `hour` : The integer hour of the pickup time. E.g., a 3:45pm taxi ride would have `15` as the hour. A 12:20am ride would have `0` .
- `day` : The day of the week with Monday=0, Sunday=6.
- `weekend` : 1 if and only if the `day` is Saturday or Sunday.
- `period` : 1 for early morning (12am-6am), 2 for daytime (6am-6pm), and 3 for night (6pm-12pm).
- `speed` : Average speed in miles per hour.

No changes are required; just run this cell.

```
In [19]: def speed(t):
             """Return a column of speeds in miles per hour."""
             return t['distance'] / t['duration'] * 60 * 60

         def augment(t):
             """Augment a dataframe t with additional columns."""
             u = t.copy()
             pickup_time = pd.to_datetime(t['pickup_datetime'])
             u.loc[:, 'hour'] = pickup_time.dt.hour
             u.loc[:, 'day'] = pickup_time.dt.weekday
             u.loc[:, 'weekend'] = (pickup_time.dt.weekday >= 5).astype(int)
             u.loc[:, 'period'] = np.digitize(pickup_time.dt.hour, [0, 6, 18])
             u.loc[:, 'speed'] = speed(t)
             return u

         train = augment(train)
         test = augment(test)
         train.iloc[0,:] # An example row
```
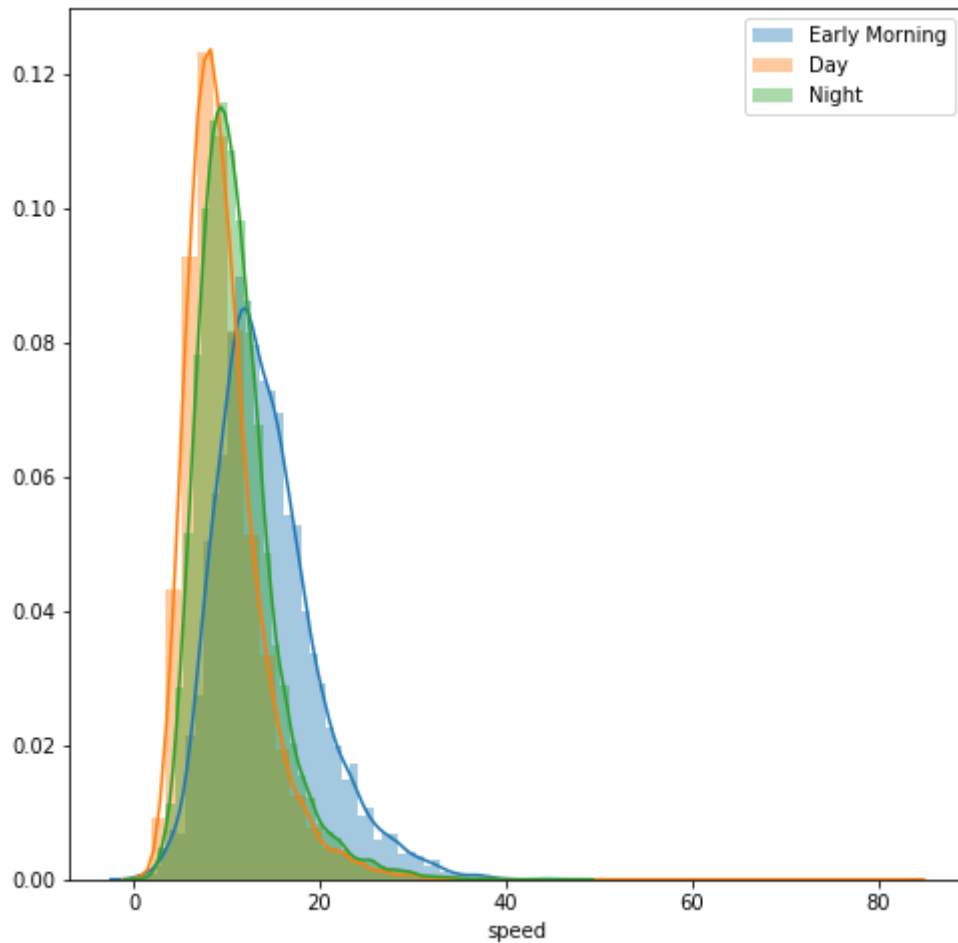
```
Out[19]: pickup_datetime      2016-01-21 18:02:20
         dropoff_datetime     2016-01-21 18:27:54
         pickup_lon                       -73.9942
         pickup_lat                         40.751
         dropoff_lon                      -73.9637
         dropoff_lat                       40.7711
         passengers                             1
         distance                            2.77
         duration                            1534
         date                          2016-01-21
         hour                                  18
         day                                    3
         weekend                                0
         period                                 3
         speed                            6.50065
         Name: 16548, dtype: object
```

## Question 3c

Use `sns.distplot` to create an overlaid histogram comparing the distribution of average speeds for taxi rides that start in the early morning (12am-6am), day (6am-6pm; 12 hours), and night (6pm-12am; 6 hours). Your plot should look like this:

```
In [20]: plt.figure(figsize=(8,8))
         morn = train[(train["hour"] >= 0) & (train["hour"] < 6)]["speed"]
         day = train[(train["hour"] >= 6) & (train["hour"] < 18)]["speed"]
         night = train[(train["hour"] >= 18) & (train["hour"] <= 23)]["speed"]
         ax = sns.distplot(morn, label="Early Morning")
         sns.distplot(day, label="Day")
         sns.distplot(night, label="Night")
         plt.legend()
         plt.show()
```



It looks like the time of day is associated with the average speed of a taxi ride.

# Question 3d

Manhattan can roughly be divided into Lower, Midtown, and Upper regions. Instead of studying a map, let's approximate by finding the first principal component of the pick-up location (latitude and longitude).

Principal component analysis (https://en.wikipedia.org/wiki/Principal_component_analysis) (PCA) is a technique that finds new axes as linear combinations of your current axes. These axes are found such that the first returned axis (the first principal component) explains the most variation in values, the 2nd the second most, etc.

Add a `region` column to `train` that categorizes each pick-up location as 0, 1, or 2 based on the value of each point's first principal component, such that an equal number of points fall into each region.

Read the documentation of `pd.qcut` (https://pandas.pydata.org/pandas-docs/version/0.23.4/generated/pandas.qcut.html), which categorizes points in a distribution into equal-frequency bins.

You don't need to add any lines to this solution. Just fill in the assignment statements to complete the implementation.

Before implementing PCA, it is important to scale and shift your values. The line with `np.linalg.svd` will return your transformation matrix, among other things. You can then use this matrix to convert points in (lat, lon) space into (PC1, PC2) space.

Hint: If you are failing the tests, try visualizing your processed data to understand what your code might be doing wrong.

*The provided tests ensure that you have answered the question correctly.*

```
In [21]: train.describe()
```

Out[21]:

|  | pickup_lon | pickup_lat | dropoff_lon | dropoff_lat | passengers | distance | |
|---|---|---|---|---|---|---|---|
| count | 53680.000000 | 53680.000000 | 53680.000000 | 53680.000000 | 53680.000000 | 53680.000000 | 5368 |
| mean | -73.981203 | 40.754681 | -73.980166 | 40.756170 | 1.662090 | 1.873336 | 66 |
| std | 0.016131 | 0.021324 | 0.017268 | 0.023868 | 1.324715 | 1.431638 | 40 |
| min | -74.017960 | 40.701378 | -74.017960 | 40.701473 | 1.000000 | 0.020000 | 6 |
| 25% | -73.992022 | 40.740189 | -73.991859 | 40.740428 | 1.000000 | 0.910000 | 36 |
| 50% | -73.982269 | 40.755581 | -73.981445 | 40.756207 | 1.000000 | 1.480000 | 57 |
| 75% | -73.970543 | 40.768520 | -73.969025 | 40.770836 | 2.000000 | 2.350000 | 86 |
| max | -73.915741 | 40.870026 | -73.916016 | 40.872025 | 6.000000 | 25.690000 | 359 |

```
In [22]: # Find the first principle component
         D = train.loc[:, ["pickup_lat", "pickup_lon"]]
         pca_n = D.count()
         pca_means = D.mean()
         X = (D - pca_means) / np.sqrt(pca_n)
         u, s, vt = np.linalg.svd(X, full_matrices=False)

         def add_region(t):
             """Add a region column to t based on vt above."""
             D = t.loc[:, ["pickup_lat", "pickup_lon"]]
             assert D.shape[0] == t.shape[0], 'You set D using the incorrect tabl
         e'
             # Always use the same data transformation used to compute vt
             X = (D - pca_means) / np.sqrt(pca_n)
             first_pc = np.dot(X, vt[0]) # select first column
             t.loc[:,'region'] = pd.qcut(first_pc, 3, labels=[0, 1, 2])

         add_region(train)
         add_region(test)
```
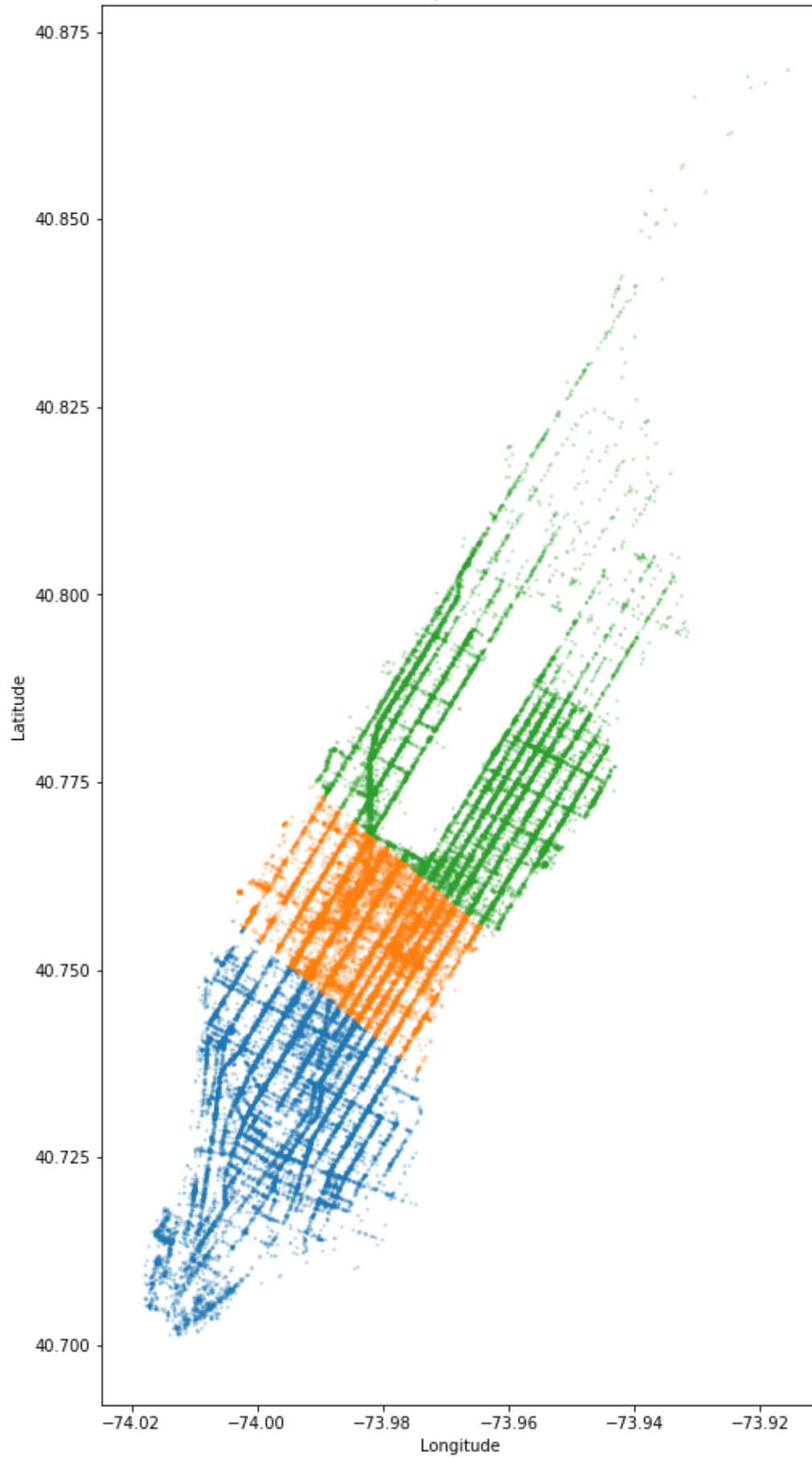
```
In [23]: grader.check("q3d")
```

Out[23]: All tests passed!

Let's see how PCA divided the trips into three groups. These regions do roughly correspond to Lower Manhattan (below 14th street), Midtown Manhattan (between 14th and the park), and Upper Manhattan (bordering Central Park). No prior knowledge of New York geography was required!

```
In [24]: plt.figure(figsize=(8, 16))
         for i in [0, 1, 2]:
             pickup_scatter(train[train['region'] == i])
```

Pickup locations

## Question 3e (ungraded)

Use `sns.distplot` to create an overlaid histogram comparing the distribution of speeds for nighttime taxi rides (6pm-12am) in the three different regions defined above. Does it appear that there is an association between region and average speed during the night?

```
In [25]: ...
Out[25]: Ellipsis
```

Finally, we create a design matrix that includes many of these features. Quantitative features are converted to standard units, while categorical features are converted to dummy variables using one-hot encoding. The `period` is not included because it is a linear combination of the `hour`. The `weekend` variable is not included because it is a linear combination of the `day`. The `speed` is not included because it was computed from the `duration`; it's impossible to know the speed without knowing the duration, given that you know the distance.

```
In [26]: from sklearn.preprocessing import StandardScaler

         num_vars = ['pickup_lon', 'pickup_lat', 'dropoff_lon', 'dropoff_lat', 'di
         stance']
         cat_vars = ['hour', 'day', 'region']

         scaler = StandardScaler()
         scaler.fit(train[num_vars])

         def design_matrix(t):
             """Create a design matrix from taxi ride dataframe t."""
             scaled = t[num_vars].copy()
             scaled.iloc[:,:] = scaler.transform(scaled) # Convert to standard uni
         ts
             categoricals = [pd.get_dummies(t[s], prefix=s, drop_first=True) for s
          in cat_vars]
             return pd.concat([scaled] + categoricals, axis=1)

         # This processes the full train set, then gives us the first item
         # Use this function to get a processed copy of the dataframe passed in
         # for training / evaluation
         design_matrix(train).iloc[0,:]
```

```
Out[26]: pickup_lon     -0.805821
         pickup_lat     -0.171761
         dropoff_lon     0.954062
         dropoff_lat     0.624203
         distance        0.626326
         hour_1          0.000000
         hour_2          0.000000
         hour_3          0.000000
         hour_4          0.000000
         hour_5          0.000000
         hour_6          0.000000
         hour_7          0.000000
         hour_8          0.000000
         hour_9          0.000000
         hour_10         0.000000
         hour_11         0.000000
         hour_12         0.000000
         hour_13         0.000000
         hour_14         0.000000
         hour_15         0.000000
         hour_16         0.000000
         hour_17         0.000000
         hour_18         1.000000
         hour_19         0.000000
         hour_20         0.000000
         hour_21         0.000000
         hour_22         0.000000
         hour_23         0.000000
         day_1           0.000000
         day_2           0.000000
         day_3           1.000000
         day_4           0.000000
         day_5           0.000000
         day_6           0.000000
         region_1        1.000000
         region_2        0.000000
         Name: 16548, dtype: float64
```

# Part 4: Model Selection

In this part, you will select a regression model to predict the duration of a taxi ride.

**Important:** *Tests in this part do not confirm that you have answered correctly. Instead, they check that you're somewhat close in order to detect major errors. It is up to you to calculate the results correctly based on the question descriptions.*

## Question 4a

Assign `constant_rmse` to the root mean squared error on the **test** set for a constant model that always predicts the mean duration of all **training set** taxi rides.

```
In [27]:  def rmse(errors):
              """Return the root mean squared error."""
              return np.sqrt(np.mean(errors ** 2))

          constant_rmse = rmse(train["duration"].mean()-test["duration"])
          constant_rmse
```

Out[27]:  399.1437572352677

```
In [28]:  grader.check("q4a")
```

Out[28]:  All tests passed!

## Question 4b

Assign `simple_rmse` to the root mean squared error on the test set for a simple linear regression model that uses only the distance of the taxi ride as a feature (and includes an intercept).

*Terminology Note*: Simple linear regression means that there is only one covariate. Multiple linear regression means that there is more than one. In either case, you can use the `LinearRegression` model from `sklearn` to fit the parameters to data.

```
In [29]:  from sklearn.linear_model import LinearRegression

          model = LinearRegression()
          simple_linreg = model.fit(train["distance"].values.reshape(-1, 1), train[
          "duration"].values.reshape(-1, 1))
          simple_rmse = rmse(simple_linreg.predict(test["distance"].values.reshape(
          -1, 1)) - test["duration"].values.reshape(-1, 1))
          simple_rmse
```

Out[29]:  276.7841105000342

```
In [30]:  grader.check("q4b")
```

Out[30]:  All tests passed!

## Question 4c

Assign `linear_rmse` to the root mean squared error on the test set for a linear regression model fitted to the training set without regularization, using the design matrix defined by the `design_matrix` function from Part 3.

*The provided tests check that you have answered the question correctly and that your `design_matrix` function is working as intended.*

```
In [31]: model = LinearRegression()
         linreg = model.fit(design_matrix(train), train["duration"])
         linear_rmse = rmse(linreg.predict(design_matrix(test)) - test["duration"
         ])
         linear_rmse
```

Out[31]: 255.19146631882776

```
In [32]: grader.check("q4c")
```

Out[32]: All tests passed!

## Question 4d

For each possible value of `period`, fit an unregularized linear regression model to the subset of the training set in that `period`. Assign `period_rmse` to the root mean squared error on the test set for a model that first chooses linear regression parameters based on the observed period of the taxi ride, then predicts the duration using those parameters. Again, fit to the training set and use the `design_matrix` function for features.

```
In [33]: model = LinearRegression()
         errors = []

         for v in np.unique(train['period']):
             v_reg = model.fit(design_matrix(train[train["period"]==v]), train[tra
         in["period"]==v]["duration"])
             errors.extend(v_reg.predict(design_matrix(test[test["period"]==v]))-t
         est[test["period"]==v]["duration"])

         period_rmse = rmse(np.array(errors))
         period_rmse
```

Out[33]: 246.62868831165173

```
In [34]: grader.check("q4d")
```

Out[34]: All tests passed!

This approach is a simple form of decision tree regression, where a different regression function is estimated for each possible choice among a collection of choices. In this case, the depth of the tree is only 1.

## Question 4e

In one or two sentences, explain how the `period` regression model above could possibly outperform linear regression when the design matrix for linear regression already includes one feature for each possible hour, which can be combined linearly to determine the `period` value.

In the original linear regression model, the line of best fit has to account for three periods at once. Since the features have different characteristics depending on period, a general model will be less accurate. By training models specifically for each period, we obtain a closer fit for each period, and by sorting data to models based on period, we are reducing the variance in the features. Both combine to improve performance.

## Question 4f

Instead of predicting duration directly, an alternative is to predict the average *speed* of the taxi ride using linear regression, then compute an estimate of the duration from the predicted speed and observed distance for each ride.

Assign `speed_rmse` to the root mean squared error in the **duration** predicted by a model that first predicts speed as a linear combination of features from the `design_matrix` function, fitted on the training set, then predicts duration from the predicted speed and observed distance.

*Hint*: Speed is in miles per hour, but duration is measured in seconds. You'll need the fact that there are 60 * 60 = 3,600 seconds in an hour.

```
In [35]: model = LinearRegression()
         speed_linreg = model.fit(design_matrix(train), train["speed"])
         speed_preds = speed_linreg.predict(design_matrix(test))
         durations = test["distance"]/(speed_preds/3600)
         speed_rmse = rmse(durations-test["duration"])
         speed_rmse
```

```
Out[35]: 243.01798368514955
```

```
In [36]: grader.check("q4f")
```

Out[36]:  All tests passed!

*Optional*: Explain why predicting speed leads to a more accurate regression model than predicting duration directly. You don't need to write this down.

My best guess: in the original model, the duration is predicted as a linear combination of the inputs. By predicting the speed then obtaining duration with the observed distance, the final prediction is further scaled with the observation itself. This makes the prediction more accurate.

## Question 4g

Finally, complete the function `tree_regression_errors` (and helper function `speed_error`) that combines the ideas from the two previous models and generalizes to multiple categorical variables.

The `tree_regression_errors` should:

- Find a different linear regression model for each possible combination of the variables in `choices`;
- Fit to the specified `outcome` (on train) and predict that `outcome` (on test) for each combination (`outcome` will be `'duration'` or `'speed'`);
- Use the specified `error_fn` (either `duration_error` or `speed_error`) to compute the error in predicted duration using the predicted outcome;
- Aggregate those errors over the whole test set and return them.

You should find that including each of `period`, `region`, and `weekend` improves prediction accuracy, and that predicting speed rather than duration leads to more accurate duration predictions.

If you're stuck, try putting print statements in the skeleton code to see what it's doing.

```
In [37]: model = LinearRegression()
         choices = ['period', 'region', 'weekend']

         def duration_error(predictions, observations):
             """Error between duration predictions (array) and observations (data
              frame)"""
             return predictions - observations['duration']

         def speed_error(predictions, observations):
             """Duration error between speed predictions and duration observation
         s"""
             return observations["distance"]/(predictions/3600) - observations['du
         ration']
             ...

         def tree_regression_errors(outcome='duration', error_fn=duration_error):
             """Return errors for all examples in test using a tree regression mod
         el."""
             errors = []
             for vs in train.groupby(choices).size().index: #values
                 v_train, v_test = train, test
                 for v, c in zip(vs, choices): #filter by choice ==
                     v_train = v_train[v_train[c]==v]
                     v_test = v_test[v_test[c]==v]
                 v_model = model.fit(design_matrix(v_train), v_train[outcome])
                 v_preds = model.predict(design_matrix(v_test))

                 errors.extend(error_fn(v_preds, v_test))
             return errors

         errors = tree_regression_errors()
         errors_via_speed = tree_regression_errors('speed', speed_error)
         tree_rmse = rmse(np.array(errors))
         tree_speed_rmse = rmse(np.array(errors_via_speed))
         print('Duration:', tree_rmse, '\nSpeed:', tree_speed_rmse)
```

```
Duration: 240.3395219270353
Speed: 226.90793945018308
```

```
In [38]: grader.check("q4g")
```
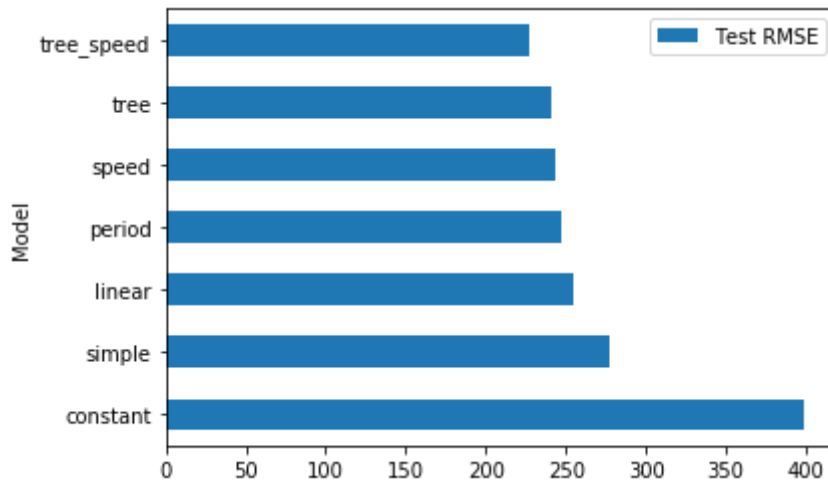
Out[38]: All tests passed!


Here's a summary of your results:

```
In [39]: models = ['constant', 'simple', 'linear', 'period', 'speed', 'tree', 'tre
         e_speed']
         pd.DataFrame.from_dict({
             'Model': models,
             'Test RMSE': [eval(m + '_rmse') for m in models]
         }).set_index('Model').plot(kind='barh');
```



# Part 5: Building on your own

In this part you'll build a regression model of your own design, with the goal of achieving even higher performance than you've seen already. You will be graded on your performance relative to others in the class, with higher performance (lower RMSE) receiving more points.

## Question 5a

In the below cell (feel free to add your own additional cells), train a regression model of your choice on the same train dataset split used above. The model can incorporate anything you've learned from the class so far.

The model you train will be used for questions 5b and 5c

```
In [40]: train_fts = design_matrix(train)
         # train_lbl = train["speed"]
         train_lbl = train["duration"]
         test_fts = design_matrix(test)
         # test_lbl = train["speed"]
         test_lbl = test["duration"]
```

```python
In [41]:  import tensorflow as tf

          from tensorflow import keras
          from tensorflow.keras import layers
          import tensorflow_docs as tfdocs
          import tensorflow_docs.plots
          import tensorflow_docs.modeling

          print(tf.__version__)
```

2.1.0

```python
In [42]:  def build_tf_model():
              model = keras.Sequential([
                  layers.Dense(64, activation='relu', input_shape=[len(train_fts.ke
          ys())]),
                  layers.Dense(64, activation='relu'),
                  layers.Dense(1)
              ])
              optimizer = tf.keras.optimizers.RMSprop(0.001)
              model.compile(loss='mse',
                            optimizer=optimizer,
                            metrics=['mae', 'mse'])
              return model

          tf_model = build_tf_model()
          tf_model.summary()
```

```
Model: "sequential"
_____
Layer (type)                 Output Shape              Param #
=================================================================
dense (Dense)                (None, 64)                2368
_____
dense_1 (Dense)              (None, 64)                4160
_____
dense_2 (Dense)              (None, 1)                 65
=================================================================
Total params: 6,593
Trainable params: 6,593
Non-trainable params: 0
_____
```
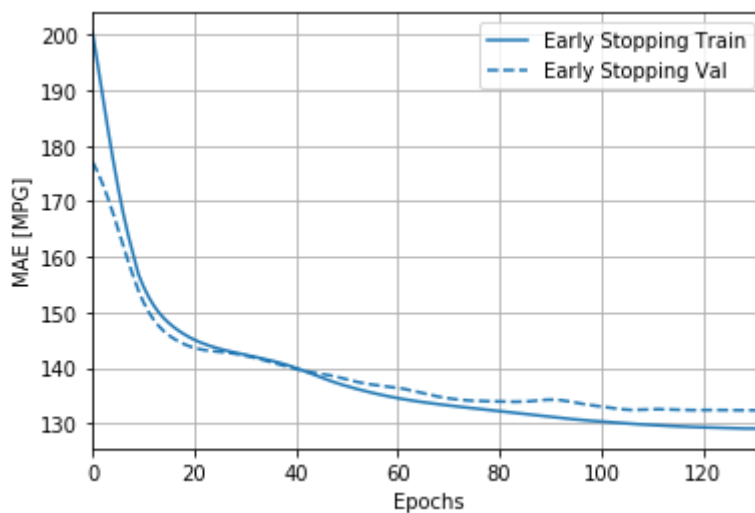
```
In [43]: tf_model = build_tf_model()
         EPOCHS = 500
         early_stop = keras.callbacks.EarlyStopping(monitor='val_loss', patience=1
         5)
         early_history = tf_model.fit(train_fts, train_lbl,
                             epochs=EPOCHS, validation_split = 0.2, verbose=0,
                             callbacks=[early_stop, tfdocs.modeling.EpochDots()])
```

```
Epoch: 0, loss:175425.4303,  mae:289.1492,  mse:175425.5312,  val_loss:
67952.9673,  val_mae:186.5232,  val_mse:67952.9688,
......................................................................

Epoch: 100, loss:35551.3135,  mae:130.4506,  mse:35551.3242,  val_loss:
36322.1962,  val_mae:135.3771,  val_mse:36322.1914,
................................
```

```
In [44]: plotter = tfdocs.plots.HistoryPlotter(smoothing_std=2)
         plotter.plot({'Early Stopping': early_history}, metric = "mae")
         plt.ylabel('MAE [MPG]')
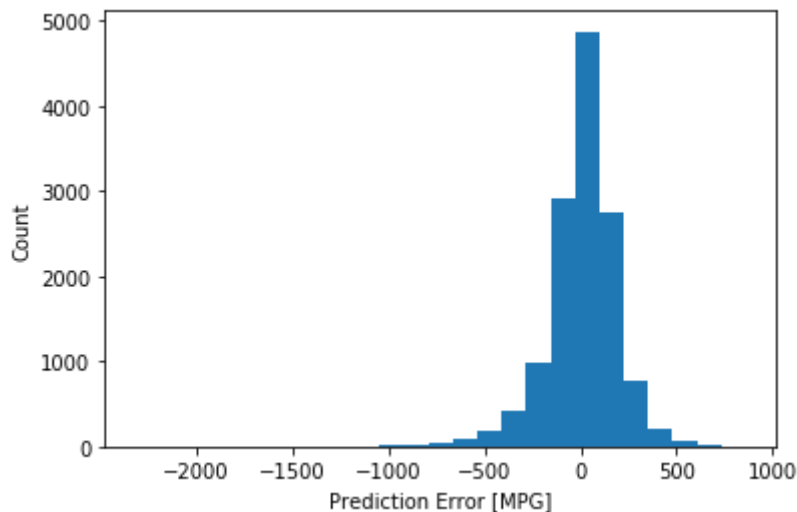```

Out[44]: Text(0, 0.5, 'MAE [MPG]')



## Question 5b

Print a summary of your model's performance. You **must** include the RMSE on the train and test sets. Do not hardcode any values or you won't receive credit.

Don't include any long lines or we won't be able to grade your response.

```
In [45]: loss, mae, mse = tf_model.evaluate(test_fts, test_lbl, verbose=2)
         print("Testing set Mean Abs Error: {:5.2f} MPG".format(mae))
         tf_train_preds = tf_model.predict(train_fts).flatten()
         tf_train_error = tf_train_preds - train["duration"]
         print("Train RMSE:", rmse(tf_train_error))
         tf_test_preds = tf_model.predict(test_fts).flatten()
         tf_test_error = tf_test_preds - test["duration"]
         print("Test RMSE:", rmse(tf_test_error))
```

```
13421/13421 - 0s - loss: 37395.0469 - mae: 133.5007 - mse: 37395.0586
Testing set Mean Abs Error: 133.50 MPG
Train RMSE: 186.00617429268308
Test RMSE: 193.37810073467472
```

```
In [46]: error = tf_test_preds - test_lbl
         plt.hist(error, bins = 25)
         plt.xlabel("Prediction Error [MPG]")
         _ = plt.ylabel("Count")
         plt.show()
```



## Question 5c

Describe why you selected the model you did and what you did to try and improve performance over the models in section 4.

Responses should be at most a few sentences

I kept the pipeline from design_matrix as I agreed with the design choices made. I thought about feature crossing latitude/longitude, but decided against it because "region" was already included and I saw no need for more fine-grained binning of locations. I then decided to use a neural network, since they can discover latent variables without manual input or feature crossing.

To this end, I used the TensorFlow library to create a neural network with two inner, Dense layers with 64 neurons each, using the relu activation function to threshold at 0 and avoid gradient vanishing. I found that the model was overfitting, so I configured the model to stop training early with a patience of 15 epochs.

(I then tried to train one neural net for each region, but found that RMSE did not improve and was actually worse even after correcting for overfitting. Thus I have kept this single-net model for submission.)

**Congratulations**! You've carried out the entire data science lifecycle for a challenging regression problem.

In Part 1 on data selection, you solved a domain-specific programming problem relevant to the analysis when choosing only those taxi rides that started and ended in Manhattan.

In Part 2 on EDA, you used the data to assess the impact of a historical event---the 2016 blizzard---and filtered the data accordingly.

In Part 3 on feature engineering, you used PCA to divide up the map of Manhattan into regions that roughly corresponded to the standard geographic description of the island.

In Part 4 on model selection, you found that using linear regression in practice can involve more than just choosing a design matrix. Tree regression made better use of categorical variables than linear regression. The domain knowledge that duration is a simple function of distance and speed allowed you to predict duration more accurately by first predicting speed.

In Part 5, you made your own model using techniques you've learned throughout the course.

Hopefully, it is apparent that all of these steps are required to reach a reliable conclusion about what inputs and model structure are helpful in predicting the duration of a taxi ride in Manhattan.

# Future Work

Here are some questions to ponder:

- The regression model would have been more accurate if we had used the date itself as a feature instead of just the day of the week. Why didn't we do that?
- Does collecting this information about every taxi ride introduce a privacy risk? The original data also included the total fare; how could someone use this information combined with an individual's credit card records to determine their location?
- Why did we treat `hour` as a categorical variable instead of a quantitative variable? Would a similar treatment be beneficial for latitude and longitude?
- Why are Google Maps estimates of ride time much more accurate than our estimates?

Here are some possible extensions to the project:

- An alternative to throwing out atypical days is to condition on a feature that makes them atypical, such as the weather or holiday calendar. How would you do that?
- Training a different linear regression model for every possible combination of categorical variables can overfit. How would you select which variables to include in a decision tree instead of just using them all?
- Your models use the observed distance as an input, but the distance is only observed after the ride is over. How could you estimate the distance from the pick-up and drop-off locations?
- How would you incorporate traffic data into the model?

```
In [ ]:   # Save your notebook first, then run this cell to generate a PDF.
          # Note, the download link will likely not work.
          # Find the pdf in the same directory as your proj3.ipynb
          grader.export("proj3.ipynb", filtering=False)
```